

UNITED STATES PATENT APPLICATION

FOR

Method and Apparatus for Fundamental Operations on Token Sequences: Computing  
Similarity, Extracting Term Values, and Searching Efficiently

INVENTOR:

Russell T. Nakano

Prepared by:

Alan Heimlich

HEIMLICH LAW

5952 DIAL WAY

SAN JOSE, CA 95129

TEL: (408) 253-3860

Docket No. 112003.P002

**Method and Apparatus for Fundamental Operations on Token Sequences: Computing  
Similarity, Extracting Term Values, and Searching Efficiently**

**RELATED APPLICATION**

**[0001]** This patent application claims priority of U.S. Provisional Application Serial No. 60/448,508 filed February 19, 2003 titled "Method and Apparatus for Operations on Token Sequences: Computing Similarity, Extracting Term Values, and Searching Efficiently", which is hereby incorporated herein by reference.

**FIELD OF THE INVENTION**

**[0001]** The present invention pertains to fundamental operations on token sequences. More particularly, the present invention relates to a method and apparatus for fundamental operations on token sequences: computing similarity, extracting term values, and searching efficiently.

**BACKGROUND OF THE INVENTION**

**[0002]** The volume of information content in our modern world is exploding at a staggering pace. In many situations, it is impossible to read, process, summarize, and extract useful meaning fast enough. This is true for text, images, sound, and video.

**[0003]** Information content is stored for two main reasons. First, content is stored for human consumption. For example, web pages, document, music, and email are stored

mainly to be redisplayed for their respective human audiences. Second, content is stored for retrieval and processing by machine. Retrieval and processing by machine requires structured information. Structured information has its data items tagged or encoded with values sufficiently uniform for machine handling. For example, XML (extensible markup language) is a kind of structured information format that is amenable for machine processing and transfer. However, most information is unstructured, and cannot be processed by machines. This presents a problem.

**[0004]** For example, consider one form of information that touches many people - email. A denizen of our increasingly interconnected world may send and receive a considerable volume of email. A tiny fraction of this email is structured with tags such as the sender, recipient(s), subject, and receipt time. The structured part of email makes it possible to sort incoming email into categories using search criteria based on the send, recipients, subject line, etc. However, little is done with the email body. The body and the attachment part of the email remain largely unstructured. As plain text, the body is searchable by wildcards, or by regular expressions, while an attachment may be scanned for known viruses. While useful, much more could potentially be done to benefit the user. For example, sorting email by the contents of the body; identifying emails that contain information of future reference value; identifying emails that have only transient value; identifying emails that require a follow-up action; and identifying emails that convey information that should be extracted, logged, and followed up. Most email users may benefit if the body could be meaningfully tagged so that a machine could perform such identifications.

**[0005]** Just as the speed and volume of email impose a daunting challenge to

individuals, similar problems afflict modern organizations around the handling of unstructured content. Companies forge agreements, establish commitments, and fulfill obligations. Legal relationships with customers, suppliers, and partners are represented in contracts. The contracts specify business terms, such as effective date, ending date, prices, quantities, and payment terms. Some contain limitation of liability clauses, while others may not. By being able to automatically scan contracts for the presence of specific types of clauses, companies may gain greater control and visibility of supplier obligations and partner commitments. However, the vast majority of contracts are stored as unstructured content. In many cases companies continue to store contracts in paper form. However, with the ability to elicit meaningful interpretations from a contract, it may become a competitive advantage to exploit machine readable contracts as a source of business virtuosity.

**[0006]** Another example is a full-text repository of issued patents which contains tags, for example for the inventor name, assignee name, and patent number. These are portions of the content that is searchable because it is structured information. The claims and description of the patent application are available in electronic form, however they are typically searchable only by Boolean search criteria. The bulk of a patent description is available as unstructured text. Someone searching must form Boolean expressions related to the presence or absence of chosen words. What is lacking is an effective approach to find explanations and claims that are "similar" to other descriptions. This presents a problem.

**[0007]** One technique for coping with unstructured textual content is text data mining.

One approach uses latent semantic analysis. This involves filtering the content to remove a list of “stop” words, such as “of,” “the,” etc. Then the remaining words are deemed to be important keywords, and their relative frequencies of appearance within a document are tallied. The numerical values map the document into a vector space, and a dimensionality reduction technique (for example, singular value decomposition (SVD)) identifies the documents that are most similar. Another approach is based on n-grams of characters and words for computing dissimilarity, equivalently, or similarity of a target text against a reference text.

**[0008]** In another approach used in the realm of plagiarism-detection, student program code is analyzed to detect plagiarism. Similarly, student term papers are analyzed with respect to subject matter, word choice, and sentence structure to identify the possibility that a given paper has partially or substantially plagiarized its content from a number of known online sources of pre-written term papers.

**[0009]** Other algorithms are used for matching images of faces. Others use an algorithm for computing similarity of an email to a spam profile.

**[0010]** The field of stylometry assigns numerical measures to attributes of one or more source texts, and performs analyses in order to gain a better understanding of such texts. For example, one application is to determine authorship of disputed works. Researchers have used principal components analysis (PCA) to analyze word frequencies of disputed texts to shed light on authorship.

**[0011]** Frequency-based approaches may not be able to analyze content at a level of detail to extract term values from textual content. This presents a problem.

## BRIEF DESCRIPTION OF THE DRAWINGS

**[0012]** The invention is illustrated by way of example and not limitation in the figures of the accompanying drawings in which:

**[0013]** Figure 1 illustrates a network environment in which the method and apparatus of the invention may be implemented;

**[0014]** Figure 2 is a block diagram of a computer system which may be used for implementing some embodiments of the invention;

**[0015]** Figure 3 illustrates eigenvalues in the complex plane for an example of the present invention;

**[0016]** Figure 4 shows a page versus page similarity comparison where exact same page to page matches (i.e. page N v. page N) have been suppressed for an example of the present invention;

**[0017]** Figure 5 illustrates a target clause with the desired business term values underlined;

**[0018]** Figure 6 illustrates one embodiment of the invention in flow chart form;

**[0019]** Figure 7, Figure 8, Figure 9, and Figure 10 illustrate embodiments of the invention showing two different possible search approaches;

**[0020]** Figure 11 through Figure 35 illustrate one embodiment of the invention as a demonstration;

**[0021]** Figure 36 through Figure 39 illustrate one embodiment of the invention showing a usage scenario;

**[0022]** Figure 40 through Figure 50 illustrate one embodiment of the invention showing

business term value extraction;

**[0023]** Figure 51 through Figure 65 illustrate one embodiment of the invention for a contract analysis;

**[0024]** Figure 66 through Figure 71 shows an overview of a procedure to compute similarity of two text blocks for one embodiment of the invention;

**[0025]** Figure 72 through 74 is an XML representation of a profile in one embodiment of the invention;

**[0026]** Figure 75 is an XML representation of a null profile in one embodiment of the invention; and

**[0027]** Figure 76 through 80 is another example of an XML representation of a profile in one embodiment of the invention.

## DETAILED DESCRIPTION

**[0028]** The present invention is a method and apparatus for fundamental operations on token sequences: computing similarity, extracting term values, and searching efficiently. Because the title of the invention is of limited length, it should not be construed as limiting the present invention to only the three operations disclosed in the title.

**[0029]** In order to explain the present invention in some of its many possible embodiments a series of explanations will be used. First, is shown how many important classes of problems can be modeled as a sequence of *tokens*. Second, is shown how to build a *probabilistic state transition model* (for example, Markov model) that uses a concept of current state that incorporates a history of recently seen tokens, which are denoted as *tuples*. Third, the probabilistic state transition model optionally uses a novel technique for assigning probabilities to the tokens that define the next state. Fourth, by combining the previous steps, a probabilistic state transition model is generated that captures the essential character of an input sequence in terms of tuples that transition to tokens. Fifth, instead of working in the tuple-token vector space naturally introduced by the probabilistic state transition model, it is shown how to generate a collection of orthogonal basis vectors that summarizes the core behavior of the input sequence - this stored mathematical representation of a sequence is denoted as its *profile*. Sixth, given a profile and a target sequence of length  $L$ , it is shown how to derive a numerical similarity score of the target relative to the profile in big-oh order  $L$  time. Seventh, given a profile it is shown how to take an input consisting of many target sequences and to efficiently search for the highest ranking similar sequences in big-oh order  $L \cdot \log(N)$  time, where  $L$  is the average length of



the sequences and N is the number of target sequences.

**[0030]** As exemplified in various embodiments of the invention, analysis of content is done at the token-sequence level. This may yield superior results compared to current approaches that are based on keyword frequencies. For example, frequency-based approaches may not be able to analyze content at a level of detail to extract term values from textual content. Term extraction uses information gleaned from a sequence of tokens to identify a target phrase. For example, find the name of a seller in the introductory clause in a purchase agreement.

**[0031]** Although this particular problem is described in terms of text examples, please observe that text merely happens to represent one instance of a general class of data type. Text is nothing more than a particular sequence of tokens. The techniques described herein can be used to analyze any other input data that can be interpreted as a sequence of tokens. For example, in this document is shown how an XML or HTML expression can be viewed as a sequence of tokens. Similarly, a word can be viewed as a sequence of characters, or a sequence of morphemes. Moving up to the document level, a purely straightforward approach would view a document as a sequence of words, but one will observe that an earlier application of this technique on words within a sentence could map to a speech acts, which in turn leads to treating a document as a sequence of speech acts.

**[0032]** Many kinds of unstructured content can be interpreted as a sequence of tokens. In the realm of text for example, any of the following can be interpreted as a token:

1. A word, capitalization left intact
2. A word, normalized uniformly to say, lower case

3. A word, stemmed or otherwise transformed to its underlying lexeme
4. A morpheme identified within an individual word
5. An individual character within a word
6. An individual character within a word, after applying one of the word-level transformations indicated above.

**[0033]** In some situations, more than one of these interpretations may be used simultaneously to improve overall accuracy and discrimination power.

## **OVERVIEW**

**[0034]** To succinctly explain the present invention, as presented here the description proceeds in three parts. First, is introduced a “profiling” module that takes input text and computes a profile, which is a combination of symbolic and numerical information that summarizes essential aspects of the input “text”. Second, is shown how a “similarity” module uses a profile to take an arbitrary input text and compute a numerical similarity score. Third, this fundamental capability is applied to solve a number of important problems in information retrieval, including recognizing similar blocks of text, extracting values of terms within a block of text, and fast indexing and retrieval of text blocks.

## **NOTA BENE**

**[0035]** Whereas other approaches may characterize an input document according to the frequencies at which predetermined keywords appear, in one embodiment the present invention uses an approach that analyzes the specific sequence of words that appear. This

may provide a substantial improvement over keyword frequency approaches because it is able to operate on smaller blocks of text while retaining the ability to discriminate between and to discern similarity among blocks of text.

**[0036]** One of skill in the art will note that the ability to discern similarity of sequences significantly enlarges the range of applicability of this technique. For example, while analysis of words is certainly one domain where sequence matters, recognizing patterns within sequences of non-verbal behaviors is important in many other domains. Fraud detection scrutinizes sequences of financial transactions. Threat detection discerns patterns within sequences of movements of people, messages, and other assets, etc. Consumer preference research studies sequences of credit purchases, movie rental decisions, web browser destinations, radio station selections, television channel choices, etc.

**[0037]** In one embodiment of the present invention is presented a technique of computing similarity of a target token sequence (for example: a document, a text block, etc.) that scales according to  $O(\log(n))$ , where  $n$  is the number of reference sequences in a repository. A logarithmic factor allows support of large repositories. For example, suppose one wants to build an automatic email query classifier, and has a training repository of  $n$  reference emails. Given an incoming email, which is called the target, the goal is to find the most similar reference email(s) in the repository to the incoming one. In the most naïve approach, one would compute the similarity of each email in the repository against the target. In this naïve approach the time to handle a query grows linearly, or  $O(n)$ . Because in one embodiment of the present invention the approach scales logarithmically, as

$O(\log(n))$ , very large repositories become practical. Alternatively, as the repository size grows one can keep response time relatively small.

## **THEORY OF OPERATION**

**[0038]** Because the present invention has many possible embodiments, Applicant has provided a brief overview of the theory of operation so that one skilled in the art will have an overview of techniques which may be used in various embodiments of the present invention. Please note that for ease of explanation, certain assumptions are made (for example, stationary processes), and that the invention is not so limited.

**[0039]** Start with a source of tokens, call it  $S$ . As has been pointed out, the tokens could be words and relevant punctuation. Alternatively, the tokens may be HTML fragments, channel selections, credit card transactions, incoming orders, troop movements, email actions, or whatever makes sense for a given problem domain. Assume that at time  $i \geq 0$ ,  $S$  emits a token  $t.i = S(i)$ . For example, the output sequence  $\{t.i ; t.i = S(i)\}$  might represent words from a block of text. Further assume that  $S$  is stationary, which means that the underlying statistics are independent of the time  $i$ . For example, if a source  $S$  generates the words of the indemnification clause from a contract, then a sentence that it generates (from the words that it emits) tend to be independent of whether it is generated early or late in the sequence.

**[0040]** Using these assumptions, one can define a Markov model that defines the probability of emitting a given token  $t$  drawn from a set  $R$ , given a prior history tuple of tokens of size  $h$ ,  $(t.0, \dots, t.h-1)$ . Call  $Q$  the set of possible tuples. In other words, one can

represent the original source S as a state machine that remembers the most recent history of h tokens emitted as a tuple  $q=(t.0,...,t.h-1)$  as a member of Q, and expresses the probability of emitting a tuple t from the set R:

$P(t, q)$  = probability of emitting token t, given the current history q.

**[0041]** Introduce a probability vector x, whose length is the number of elements of Q, and where the j-th element  $x.j$  represents the probability of being in the state of the j-th tuple. Because we assume the source S to be stationary the probability transition function P can be represented by a real matrix A, where the (j,k)-th element of the matrix is,

$A(j, k)$  = probability that S emits token r.k, given that the current state is given by tuple q.j.

**[0042]** In this formulation, one can see that,

$$\text{sum}(k; A(j, k)) = 1$$

since the token that the source emits must be one of the tokens contained in the set R. That is the typical assumption for a Markov model. As noted later, this assumption is not essential, because in some situations (embodiments of the invention) one can truncate the matrix A by deleting rows and/or columns that are deemed non-essential. Another property of a Markov model is  $0 \leq A(j, k) \leq 1$ , and for now this remains so in our model, although we do not particularly require it.

**[0043]** The matrix A has the number of rows equal to the number of possible tuples, and the number of columns equal to the number of possible tokens. The matrix A maps a probability vector x in tuple space, where each element  $x_j$  represents the probability of the current state being the j-th tuple, to a probability vector y in token space, where  $y_k$  represents the probability of the next token being the k-th token. This can be expressed,

$$x \cdot \text{transpose } A = y \cdot \text{transpose}$$

where x is the probability vector in tuple space, and y is the probability vector in token space.

**[0044]** Since the number of tuples exceeds the number of tokens, in the typical case A will have more rows than columns. If A isn't square, then solving the singular value decomposition (SVD) problem for A yields the factorization,

$$A = U D V^{\text{hermitian}}$$

where D is diagonal, U and V are orthogonal (unitary), and  $V^{\text{hermitian}}$  is formed by transposing V and taking the complex conjugate of its elements.

**[0045]** One will observe that A will usually be very sparse, and by renumbering the tuples it is possible to place the tuples that occur most frequently in the low numbered rows. If it also occurs that some tokens are far less frequent than others, then by a similar

renumbering of the token assignments to the columns of A one can congregate the non-zero elements of A into the low-numbered columns. Applicant has observed that after renumbering, that A can be reduced to a square matrix by dropping some number of high-numbered rows and columns. Applicant's experience has been that as long as renumbering places the high-valued, non-zero elements in the low-numbered rows and columns, then similarity results tend to be fairly robust to deleting high-numbered rows and columns.

**[0046]** If by this procedure one creates a square matrix A that approximates the original A, one can solve the nonsymmetric eigenproblem (NEP). The eigenproblem also yields a factorization,

$$A = U D V.\text{hermitian}$$

where D is diagonal with eigenvalues down the diagonal, U and V are orthonormal.

**[0047]** One will observe that either solving the SVD problem for a non-square A, or solving the NEP for a square A, both formulations provide factorizations with orthogonal U and V. Regardless of which technique is used to solve the problem, the key result is that U and V are orthogonal. We use this property as follows. Premultiply the expression above by U.hermitian and postmultiply by V, and the result is

$$U.\text{hermitian} A V = U.\text{hermitian} U D V.\text{hermitian} V$$

$$= (U.\text{hermitian } U) D (V.\text{hermitian } V)$$

$$= D$$

where one may use the fact that  $U.\text{hermitian } U = I$  and  $V.\text{hermitian } V = I$ , where  $I$  is the identify matrix.

**[0048]** Now one can diagonalize  $A$  by the use of  $U$  and  $V$ . Whether one uses SVD or NEP is determined more by the pragmatic decisions of whether one truncates  $A$  to make it square, or whether one prefers to deal with a non-square  $A$ . The other reason for not favoring one technique over the other is that it turns out that neither the eigenvalues (NEP) or the singular values (SVD) come into play later, when we compute the similarity measure. We are interested in  $U$  and  $V$  because they diagonalize  $A$ .

**[0049]** Diagonalizing a matrix has a useful interpretation. Multiplying by  $U$  and  $V$  have the effect of changing to a new basis, where under the new basis all the “energy” of the matrix has been channeled into independent dimensions.

**[0050]** A key insight from this approach is that  $U$  and  $V$  provide a custom-constructed basis with which to describe the essential character of the original  $A$ . That is why we call  $U$  and  $V$  the “profile” of  $A$ . In a sense,  $U$  and  $V$  represent an substantially optimal basis for the transformation underlying the matrix  $A$ .

**[0051]** Now is introduced the concept of similarity.  $U$ - $V$  factorization and interpretation has been previously discussed because that will assist one to understand what similarity is. Suppose one has another matrix  $B$ , which arises from the probabilistic transformation model from another source of tokens. For example, this could be derived from the



sequence of words in a paragraph of an unknown contract. Compute the similarity of the transformation represented by B by using the profile computed from A, namely the U and V matrices. Specifically, compute the matrix,

$$G = U \cdot \text{hermitian } B \cdot V$$

**[0052]** One knows that if B were identical to A, then U and V would diagonalize B. In other words, all the energy of B would also be forced down the diagonal of G. If B were substantially similar to A, then one may expect that the result G would be nearly diagonal. Applicant's empirical results support this.

**[0053]** This is the rationale behind the similarity measure, defined as,

$$\text{Similarity} = \sum(i; || u.i.\text{conjugate} * b * v.i || ^2)$$

**[0054]** The i-th term in the sum is precisely the i-th diagonal in the matrix  $G = U \cdot \text{hermitian } B \cdot V$ . If A is square, then i ranges from 1 to the size of A. Otherwise i goes up to number of tokens or the number of tuples, whichever is smaller.

**[0055]** So, recapping the overall procedure. Observe a sequence of tokens from the set R, and form the history tuples from the set Q. If it is assumed that the underlying process that generated the tokens is stationary, then one can compute the corresponding A matrix. Next, factor matrix A into the orthogonal matrices U and V. Taken together, R, Q, U and V define the profile of our original sequence. Armed with this profile one can proceed to observe another sequence, from which one may compute its corresponding transition matrix B according to the tokens R and tuples Q. Compute similarity between the first

sequence and the second sequence as described earlier. In doing so, one is computing a degree to which U and V diagonalize B. This is the "essence" of similarity.

**[0056]** One of skill in the art will also note that this analysis also tells one that the presence of higher magnitude off-diagonal elements in G, which we might call the "similarity matrix," indicates a greater deviation from the original A matrix.

## **PROFILE**

**[0057]** In one embodiment of the invention, this is a procedure to compute a profile.

1. Start with an input block of text.
2. Tokenize the text so that words and sentences are defined. Introduce a distinguished token to mark the end of sentence. In the following discussion, the token TERM is used to indicate the end of a sentence. For illustration purposes, convert words in the block of text into lowercase. Any token that contains a capital letter can safely be used as a special token. The result of this step is to represent the input block of text as a sequence of tokens  $t.i$ , where some of the tokens are sentence delimiters:  $\{t.i : i=0, \dots, n\}$ . (Note that some form of stemming that maps different forms of the same word or verb phrase on to a common root improves the ability of this approach to detect patterns.)
3. Choose a small number  $h>0$ , which shall be called our history window. For English text, it has been observed that  $h=2$  works well.
4. Each token  $t.i$  for  $i>h$ , has  $h$  predecessors  $q.i = (t.i-h, t.i-h+1, \dots, t.i-1)$ . For completeness, define

q.0 to be (TERM, ... TERM).

q.1 to be (TERM, ..., t0)

q.h to be (t.0, ..., t.h-1)

5. Using this approach, for every position t.i in the input sequence  $i=0, \dots, n$ , there is a predecessor history q.i. Refer to the i-th predecessor p.i as the i-th tuple, tuple because it is an h-tuple of tokens that precede token t.i.
6. For each token t.i for  $i \geq 0$ , there is a next token in the sequence  $r.i = t.i$ , which is the token itself.
7. Therefore, each token t.i has a predecessor q.i, defined above, and a subsequent token r.i.
8. The collection of tuples,  $Q = \{q.i : i=0, \dots, n\}$ , defines a set of states in the tuple space. Similarly, the collection of tokens,  $R = \{r.i : i=0, \dots, n\}$ , defines a set of states in the token space. Each token t.i represents a transition from the predecessor h-history of tokens, to the current token.
9. Define a stationary Markov model that has a probability of transitioning from a state q.i contained in the set Q, to an element in the token set R:

$P(q, r)$  = probability of that the next token is r, given that the previous tokens are described by the tuple q.

10. The probability model is defined over the cross product of the tuples  $Q$  and the tokens  $R$ . In practice, there will be a subset of the observed tuples  $Q$  that are the most frequently seen tuples and a similar most-frequent subset of the tokens. Call these reduced sets  $Q.star$  and  $R.star$ . The primary consideration in choosing a strict subset  $Q.star$  of  $Q$ , or subset  $R.star$  of  $R$ , is to keep the computational load manageable. As compute power increases, or as the time available to perform computations increases, or as sparse matrix techniques becomes more widely available, then the need to reduce the tuple and token sets becomes less pressing. (Note that for certain applications some tokens may be deemed highly important even if they don't fall within the most frequent tuple set  $Q.star$  or the most frequent token set  $R.star$ . In other embodiments of the invention the sets  $Q.star$  and  $R.start$  may be augmented with such important tuples or tokens.)
11. Suppose the  $N$  is the reduced size of  $Q.star$  and  $R.star$ . For convenience label the tuples and tokens according to an indexing scheme. In practice it works well to designate item 0 to correspond to the most frequent tuple and token. Form the matrix  $a[i,j]$ ,
- where  $a[i,j]$  = probability of transitioning from tuple.i to token.j.

12. Observe that  $a[]$  is an  $N \times N$  matrix. In practice, the transition matrix  $a[]$  tends to contain many zeros. In other words, the matrix is sparse.
13. Perform a principal component analysis on the matrix  $a[]$ . This is equivalent to computing the eigenvalues  $w.i$ , left eigenvectors  $u.i$ , and right eigenvectors  $v.i$  of the matrix  $a$ . In practice, it has been seen that the eigenvalues and eigenvectors are complex. It is convenient to normalize the eigenvectors to be of unit length. By virtue of the PCA and/or eigensystem analysis ( i.e. solving this eigenproblem), the eigenvectors will be orthogonal. In practice, there will be  $C \leq N$  eigenvector pairs. Since the matrix  $a$  is real, the eigenvalues will either be real, in which case there is a corresponding left and right eigenvector pair, or there will be a complex conjugate pair of eigenvalues and a corresponding pair of complex left and right eigenvectors.
14. Define the "profile" of the input text block as follows: the  $N$  tuples and tokens chosen above, and the  $C$  eigenvector pairs computed above.  $C$  which denotes a component contains an eigenvalue and corresponding left and right eigenvectors. Thus, an  $i$ th component is  $(w.i, u.i, v.i)$ .

**[0058]** Using the profile approach discussed above, consider this example.

Consider the famous opening from a Dr. Seuss classic:

"I can read in red. I can read in blue. I can read in pickle  
color too."

After tokenizing, the following sequence of tokens is obtained, including the special token TERM that denotes the sentence delimiter.

i can read in red TERM i can read in blue TERM i can read in  
pickle color too TERM

For this example, a history window,  $h = 2$  was used. Construct the predecessor tuples  $q_i$  and tokens  $t_i$  as shown below in Table 1.

**Table 1. Tuple to token state transitions.**

i	Predecessor (p.i)	Token (t.i)
0	TERM TERM	i
1	TERM i	can
2	i can	read
3	can read	in
4	read in	red
5	in red	TERM
6	red TERM	i
7	TERM i	can
8	i can	read
9	can read	in
10	read in	blue
11	in blue	TERM
12	blue TERM	i
13	TERM i	can
14	i can	read
15	can read	in
16	read in	pickle
17	in pickle	color
18	pickle color	too
19	color too	TERM

The set of tokens R can be written as shown in Table 2.

**Table 2. The set of tokens and frequency of occurrence.**

j	Tokens (r.j)	Frequency
0	i	3
1	can	3
2	read	3
3	in	3
4	TERM	3
5	red	1
6	blue	1
7	pickle	1
8	color	1
9	too	1

The set of tuples Q can be written as shown in Table 3.

**Table 3. The set of tuples and their frequency of occurrence.**

j	Tuples (q.i)	Frequency
0	TERM i	3
1	i can	3
2	can read	3
3	read in	3
4	in red	1
5	red TERM	1
6	in blue	1
7	blue TERM	1
8	in pickle	1
9	pickle color	1
10	color too	1
11	TERM TERM	1

**[0059]** Define a stationary Markov model that transitions from state q.i to r.j with probability  $P_{i,j}$ .

For example, the probability  $P_{3.5}$  of transitioning from tuple  $q_{.3}$ , “read in,” (see  $j=3$  in Table 3) to token  $r_{.5}$ , “red,” (see  $j=5$  in Table 2) is  $1/3$ . The  $1/3$  is derived as follows. From Table 1 we see that at  $i=4$ , 10, and 16 is the  $p.i$  term “read in”. Similarly, at these three values of  $i=4$ , 10, and 16, only at  $i=4$  is the  $t.i$  term “red”. At  $i=10$  it is “blue” and at  $i=16$  it is “pickle”. Thus the probability of transition from “read in” to “red” is  $1/3$ .

As another example, the probability  $P_{3.4}$  of transitioning from “read in,” to “TERM,” is 0. Again looking at Table 1 “read in” never transitions to TERM (i.e.  $0/3$ ).

**Table 4. The transition probability matrix.**

			Tokens r.j									
			j=0	1	2	3	4	5	6	7	8	9
			i	can	read	in	TERM	red	blue	pickle	color	too
T u p l e s  (q.i)	i=0	TERM i		1								
	1	i can			1							
	2	can read				1						
	3	read in					0.33	0.33	0.33			
	4	in red					1					
	5	red TERM	1									
	6	in blue					1					
	7	blue TERM	1									
	8	in pickle								1		
	9	pickle color									1	
	10	color too					1					
	11	TERM TERM										

One of skill in the art will recognize that in the absence of truncating terms that the summation across any row in the transition probability matrix (see for example Table 4) will be 1.

**[0060]** Generally, longer input sequences will have many more tuples than tokens,



because a tuple derives from several tokens in the “history window,”  $h$ . In its full form the transition matrix has a number of rows equal to the number of tuples. Although this example is relatively small, the technique of reducing the computational load will be illustrated. One can approximate the 12x10 transition matrix by a 10x10 transition matrix. In practice, one may find that the results do not change much, especially if the most frequent tuples and tokens are retained. One such approximation with smaller dimension is the following as illustrated in Table 5.

**Table 5. A probability transition matrix of reduced dimensionality.**

	i	can	read	in	TERM	red	blue	pickle	color	too
TERM i		1								
i can			1							
can read				1						
read in						0.33	0.33	0.33		
in red					1					
in blue					1					
in pickle									1	
pickle color										1
color too					1					
TERM TERM										

**[0061]** Solving the eigenproblem arising from this matrix yields 10 complex-valued eigenvalues, as shown in Table 6. The location of the eigenvalues in the complex plane are shown in Figure 3.

**Table 6. Eigenvalues.**

Component	Real Part	Imaginary Part
0	-0.5	0.87
1	-0.5	-0.87
2	1	0
3	0	0
4	1	0
5	-1	0
6	0	0
7	1	0
8	1	0
9	1	0

**[0062]** For convenience, one can normalize all the eigenvectors to unit magnitude.

Each eigenvalue has a corresponding left eigenvector and right eigenvector. The components of the left eigenvector correspond to the tuples. Similarly, the components of the right eigenvector correspond to the tokens. Together, the left and right eigenvectors corresponding to each eigenvalue can be interpreted in terms of the tuple-to-token state transitions. This is especially true in a simple example such as this one. For example, it can be seen that the unit eigenvalue  $1+0i$  has several left-right eigenvector pairs. One of them has a left eigenvector with unit magnitude corresponding to the "i can" tuple, and a right eigenvector has unit magnitude corresponding to the "read" token.

**Table 7. An eigenvalue with its left and right eigenvalues.**

**Eigenvalue**

Component	Real part	Imaginary part
7	1	0

**Left eigenvector**

	Re	Im
I can	1	0

**Right eigenvector**

	Re	Im
read	1	0

This arises from the original transition matrix due to the “i can read” token sequence.

Another eigenvector pair associated with the  $1+0i$  eigenvalue has the left eigenvector with unit magnitude in the “pickle color” direction, with corresponding right eigenvector with unit magnitude in the “too” tuple direction. Just as in the “I can read” example, this component of the eigenproblem arises from the “pickle color too” sequence.

**Table 8. Another eigenvalue, with left and right eigenvectors.**

**Eigenvalue**

Component	Real part	Imaginary part
8	1	0

**Left eigenvector**

	Re	Im
pickle color	1	0

**Right eigenvector**

	Re	Im
too	1	0

**[0063]** An interesting pair of left and right eigenvectors is shown below in Table 9.

**Table 9. A complex eigenvalue, with left and right eigenvectors.**

**Eigenvalue**

Component	Real part	Imaginary part
0	-0.5	0.866

**Left eigenvector**

	Re	Im
read in	-0.655	0

TERM i	0.327	-0.566
in red	0.109	0.189
color too	0.109	0.189
in blue	0.109	0.189

#### Right eigenvector

	Re	Im
can	-0.447	0
blue	0.224	0.387
pickle	0.224	0.387
red	0.224	0.387
TERM	0.224	-0.387

This example has illustrated how to construct a "profile" of the source text in terms of tuples, tokens, eigenvalues, left eigenvectors, and right eigenvectors. Computing a profile entails solving the eigenproblem of the state transition matrix, and may be computationally intensive. However, observe that for most inputs the transition matrix tends to have many zero elements. That is, the matrix is sparse, which offers many opportunities for computationally efficient algorithms.

**[0064]** In the next section, is a discussion on how to use a profile to compute a numerical measure of how "similar" another token sequence is to the one used to construct the profile. The similarity computation does not require solving an eigenproblem. As a consequence, it is relatively inexpensive from a computational standpoint. This may be advantageous in many applications where a profile is computed once and used many times to measure similarity against other token sequences.

## SIMILARITY

**[0065]** Shown here is how in one embodiment of the invention, an input block of text, which one may refer to as the "target," may be used to compute a similarity score against a

given profile.

1. Process the target block of text in the same manner as was done in the process to compute the profile, up to the point where the transition probability matrix was derived. To avoid confusion, refer to this transition probability matrix as  $b()$ ,

$b(p, q)$  = observed probability of tuple  $p$  transitioning to token  $q$  in the target block of text.

2. Obtain a previously computed profile, consisting of the tuples, tokens, eigenvalues, left eigenvectors, and right eigenvectors. Refer to the text used to create the profile as the reference text.
3. Compute the similarity score of the target text against the reference as follows. For each eigenvalue  $w_i$  in the profile, obtain the left eigenvector  $u_i$  and right eigenvector  $v_i$ .

Similarity score =  $\sum_i (|| u_i \cdot \text{conjugate} * a * v_i ||^2)$ ,

where each term of the summation is the transition matrix premultiplied by the complex conjugate of the  $i$ -th left eigenvector, and postmultiplied by the  $i$ -th right

eigenvector. The norm squared  $||.||^2$  is the square of the magnitude of the complex number resulting from the  $i$ -th term in the sum.

4. A low similarity score indicates that the target text has little or nothing in common with the reference text. A high similarity score indicates that there are many common tuple-token combinations that appear in the target text, that also appear in the reference text.

**[0066]** Here is an example of how in one embodiment of the invention, the similarity computation works. Consider the text of the second page of the Dr. Seuss work:

I can read in bed. And in purple. And in brown. I can read  
in a circle and upside down!

As before, tokenize and form the transition matrix for this sequence. Observe that the matrix is sparse, which lends itself to computationally efficient routines.

**Table 10. Probability transition matrix of the target text.**

	can	read	in	bed	a	TERM	purple	circle	brown	circle	...
TERM i	1										
i can		1									
can read			1								
read in				0.5	0.5						
TERM and			1								
in purple						1					
in brown						1					
in bed						1					

and in							0.5		0.5		
in a								1			
...											

The similarity score may be obtained by using the formula,

$$\text{Similarity score} = \sum(i, || u.i.\text{conjugate} * a * v.i ||^2).$$

The i-th term of the sum uses the i-th left eigenvector of the profile to premultiply the transition matrix, and the i-th right eigenvector of the profile to postmultiply the transition matrix. For example, component #7 of the profile has the following left and right eigenvectors as shown in Table 11:

**Table 11. Complex conjugate of left eigenvector, and right eigenvector.**

Left eigenvector (complex conjugate)

	Re	Im
I can	1	0

Right eigenvector

	Re	Im
read	1	0

Observe that the transition matrix has value  $a[(i \text{ can}), \text{read}] = 1$ , hence that component contributes 1 to the sum, as the squared magnitude. In general,

$$|| x + i y ||^2 = x^2 + y^2$$

where x is a real part, and y is an imaginary part. In this notation, "i" denotes complex math and not an index. For example,  $1 + 0i$  has a squared magnitude of 1.

So  $1 + 0i$  has a squared magnitude of 1.

In contrast, component #8 contributes nothing to the similarity score because neither the tuple (pickle color), nor the token “too” appears in the current sample.

Component #0 premultiplies by the following left eigenvector as shown in Table 12.

**Table 12. Left eigenvector corresponding to a complex eigenvalue.**

Left eigenvector (complex conjugate)

	Re	Im
read in	-0.655	0
TERM i	0.327	0.566
in red	0.109	-0.189
color too	0.109	-0.189
in blue	0.109	-0.189

This yields the following row vector as shown in Table 13.

**Table 13. Result of left eigenvector and transition probability matrix.**

	can	read	in	bed	a	TERM	purple	circle	brown	circle	...
Re	0.327			0.327	0.327						
Im	0.566			0	0						

Multiply by the right eigenvector as shown in Table 14.

**Table 14. Right eigenvector.**

Right eigenvector

	Re	Im
can	-0.447	0
blue	0.224	0.387
pickle	0.224	0.387
red	0.224	0.387
TERM	0.224	-0.387



The sole non-zero term resulting from the row vector and the column vector arises from the “can” term,  $(0.327 + 0.566i)*(-0.447 + 0i) = 0.146 + 0.253i$ . The contribution of the squared magnitude is 0.085.

Combining all the components of the profile, the similarity score for page 1 relative to the computed profile for page 0 is 2.240. In practice it has been observed that a similarity score greater than one tends to indicate a noticeable congruence between the profile text and a given target.

In the illustration below are the compute profiles for all the pages of this venerable Dr. Seuss classic, and also are computed the similarity of all the other pages in the work.

**Table 15. Text of “I can read with my eyes shut!”**

Page	Text
0	I can read in red. I can read in blue. I can read in pickle color too.
1	I can read in bed. And in purple. And in brown. I can read in a circle and upside down!
2	I can read with my left eye. I can read with my right. I can read Mississippi with my eyes shut tight.
3	Mississippi, Indianapolis and hallelujah, too. I can read them with my eyes shut! That is very hard to do!
4	But it's bad for my hat and makes my eyebrows get red hot. So reading with my eyes shut I don't do an awful lot.
5	And when I keep them open I can read with much more speed. You have to be a speedy reader 'cause there's so, so much to read!
6	You can read about trees and bees and knees. And knees on trees! And bees on threes!
7	You can read about anchors. And all about ants. You can read about ankles! And crocodile pants.
8	You can read about hoses and how to smell roses and what you should do about owls on noses.

9	Young cat! If you keep your eyes open enough, oh, the stuff you will learn! The most wonderful stuff! You'll learn about fishbones and wishbones. You'll learn about trombones, too. You'll learn about Jake the Pillow Snake and all about Foo-Foo the Snoo.
10	You'll learn about ice. You can learn about mice. Mice on ice. And ice on mice.
11	You can learn about the price of ice.
12	You can learn about sad and glad and mad. There are so many things you can learn about. But you'll miss the best things if you keep your eyes shut.
13	The more that you read the more things you will know. The more that you learn, the more places you'll go.
14	You might learn a way to earn a few dollars. Or how to make doughnuts or kangaroo collars.
15	You can learn to read and play a Hut-Zut if you keep your eyes open. But not with them shut.
16	If you read with your eyes shut you're likely to find that the place where you're going is far, far behind.
17	So that's why I tell you to keep your eyes wide. Keep them wide open at least on one side.

**[0067]** The results are shown in Figure 4, where the base of the graph depicts all the page combinations, and the height of the bar represents the similarity score computed using the procedure above. Note that Figure 4 shows a page versus page similarity comparison where exact same page to page matches (i.e. page N v. page N) have been suppressed.

**[0068]** Some of the notable matches from the above example are shown below in Table 16, arranged with the highest similarity pairs at the top. One can observe that as the similarity measure diminishes the profile text and the target text become less alike.

**Table 16. Similarity between profile and a target.**

Profile Text	Target Text	Similarity
i can read in red. i can read in blue. i can read in pickle color too.	i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	2.24
i can read in red. i can read in blue. i can read in pickle color too.	i can read with my left eye. i can read with my right. i can read mississippi with my eyes shut tight.	1.56

you can learn about the price of ice.	you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	1.31
you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	young cat! if you keep your eyes open enough, oh, the stuff you will learn! the most wonderful stuff! you'll learn about fishbones and wishbones. you'll learn about trombones, too. you'll learn about jake the pillow snake and all about foo-foo the snoo.	1.30
you can read about anchors. and all about ants. you can read about ankles! and crocodile pants.	you can read about trees and bees and knees. and knees on trees! and bees on threes!	0.95
i can read with my left eye. i can read with my right. i can read mississippi with my eyes shut tight.	mississippi, indianapolis and hallelujah, too. i can read them with my eyes shut! that is very hard to do!	0.94
i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	i can read in red. i can read in blue. i can read in pickle color too.	0.81
you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	0.60
mississippi, indianapolis and hallelujah, too. i can read them with my eyes shut! that is very hard to do!	i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	0.58
you can read about trees and bees and knees. and knees on trees! and bees on threes!	you can read about anchors. and all about ants. you can read about ankles! and crocodile pants.	0.56
you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	you can learn about the price of ice.	0.56
if you read with your eyes shut you're likely to find that the place where you're going is far, far behind.	you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	0.50
you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	you can learn about the price of ice.	0.49
you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	0.45
i can read with my left eye. i can read with my right. i can read mississippi with my eyes shut tight.	i can read in red. i can read in blue. i can read in pickle color too.	0.38

young cat! if you keep your eyes open enough, oh, the stuff you will learn! the most wonderful stuff! you'll learn about fishbones and wishbones. you'll learn about trombones, too. you'll learn about jake the pillow snake and all about foo-foo the snoo.	you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	0.37
i can read with my left eye. i can read with my right. i can read mississippi with my eyes shut tight.	i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	0.28
you can read about hoses and how to smell roses and what you should do about owls on noses.	you can learn about the price of ice.	0.28
i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	mississippi, indianapolis and hallelujah, too. i can read them with my eyes shut! that is very hard to do!	0.26
you can read about hoses and how to smell roses and what you should do about owls on noses.	you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	0.26
young cat! if you keep your eyes open enough, oh, the stuff you will learn! the most wonderful stuff! you'll learn about fishbones and wishbones. you'll learn about trombones, too. you'll learn about jake the pillow snake and all about foo-foo the snoo.	you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	0.20
you can read about anchors. and all about ants. you can read about ankles! and crocodile pants.	you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	0.14
you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	you can read about trees and bees and knees. and knees on trees! and bees on threes!	0.13
you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	you can read about hoses and how to smell roses and what you should do about owls on noses.	0.13
you can read about trees and bees and knees. and knees on trees! and bees on threes!	you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	0.12
but it's bad for my hat and makes my eyebrows get red hot. so reading with my eyes shut i don't do an awful lot.	you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	0.11

you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	the more that you read the more things you will know. the more that you learn, the more places you'll go.	0.09
so that's why i tell you to keep your eyes wide. keep them wide open at least on one side.	you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	0.08
i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	and when i keep them open i can read with much more speed. you have to be a speedy reader 'cause there's so, so much to read!	0.08
you can read about hoses and how to smell roses and what you should do about owls on noses.	you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	0.08
and when i keep them open i can read with much more speed. you have to be a speedy reader 'cause there's so, so much to read!	i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	0.07
so that's why i tell you to keep your eyes wide. keep them wide open at least on one side.	you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	0.07
you can read about trees and bees and knees. and knees on trees! and bees on threes!	you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	0.06
you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	you can read about hoses and how to smell roses and what you should do about owls on noses.	0.06
you might learn a way to earn a few dollars. or how to make doughnuts or kangaroo collars.	you can learn about the price of ice.	0.05
you might learn a way to earn a few dollars. or how to make doughnuts or kangaroo collars.	you can read about hoses and how to smell roses and what you should do about owls on noses.	0.05
you can learn about the price of ice.	if you read with your eyes shut you're likely to find that the place where you're going is far, far behind.	0.04
you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	you can read about trees and bees and knees. and knees on trees! and bees on threes!	0.04
young cat! if you keep your eyes open enough, oh, the stuff you will learn! the most wonderful stuff! you'll learn about fishbones and wishbones. you'll learn about trombones, too. you'll learn about jake the pillow snake and all about foo-foo the snoo.	the more that you read the more things you will know. the more that you learn, the more places you'll go.	0.04

you can learn about the price of ice.	you might learn a way to earn a few dollars. or how to make doughnuts or kangaroo collars.	0.04
you can learn about the price of ice.	young cat! if you keep your eyes open enough, oh, the stuff you will learn! the most wonderful stuff! you'll learn about fishbones and wishbones. you'll learn about trombones, too. you'll learn about jake the pillow snake and all about foo-foo the snoo.	0.03
you can read about anchors. and all about ants. you can read about ankles! and crocodile pants.	and when i keep them open i can read with much more speed. you have to be a speedy reader 'cause there's so, so much to read!	0.03
you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	if you read with your eyes shut you're likely to find that the place where you're going is far, far behind.	0.02
you might learn a way to earn a few dollars. or how to make doughnuts or kangaroo collars.	you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	0.02
but it's bad for my hat and makes my eyebrows get red hot. so reading with my eyes shut i don't do an awful lot.	and when i keep them open i can read with much more speed. you have to be a speedy reader 'cause there's so, so much to read!	0.02
young cat! if you keep your eyes open enough, oh, the stuff you will learn! the most wonderful stuff! you'll learn about fishbones and wishbones. you'll learn about trombones, too. you'll learn about jake the pillow snake and all about foo-foo the snoo.	you can learn about the price of ice.	0.02
mississippi, indianapolis and hallelujah, too. i can read them with my eyes shut! that is very hard to do!	you can learn about the price of ice.	0.02
you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	if you read with your eyes shut you're likely to find that the place where you're going is far, far behind.	0.02
you can learn to read and play a hut-zut if you keep your eyes open. but not with them shut.	you might learn a way to earn a few dollars. or how to make doughnuts or kangaroo collars.	0.01
you can read about anchors. and all about ants. you can read about ankles! and crocodile pants.	but it's bad for my hat and makes my eyebrows get red hot. so reading with my eyes shut i don't do an awful lot.	0.01

you can read about hoses and how to smell roses and what you should do about owls on noses.	young cat! if you keep your eyes open enough, oh, the stuff you will learn! the most wonderful stuff! you'll learn about fishbones and wishbones. you'll learn about trombones, too. you'll learn about jake the pillow snake and all about foo-foo the snoo.	0.01
you might learn a way to earn a few dollars. or how to make doughnuts or kangaroo collars.	you can read about trees and bees and knees. and knees on trees! and bees on threes!	0.01
you might learn a way to earn a few dollars. or how to make doughnuts or kangaroo collars.	you can learn about sad and glad and mad. there are so many things you can learn about. but you'll miss the best things if you keep your eyes shut.	0.01
the more that you read the more things you will know. the more that you learn, the more places you'll go.	you can learn about the price of ice.	0.01
you'll learn about ice. you can learn about mice. mice on ice. and ice on mice.	so that's why i tell you to keep your eyes wide. keep them wide open at least on one side.	0.01
you can learn about the price of ice.	mississippi, indianapolis and hallelujah, too. i can read them with my eyes shut! that is very hard to do!	0.01
you can read about anchors. and all about ants. you can read about ankles! and crocodile pants.	i can read in bed. and in purple. and in brown. i can read in a circle and upside down!	0.01

## COMPUTATIONAL COST

**[0069]** An estimate concerning the number operations required to compute a similarity measure is now discussed. Suppose that one has already computed a profile, and has decided to retain  $k$  components in the profile. Recall that this means that these are stored:  $k$  eigenvalues,  $k$  left eigenvectors, and  $k$  right eigenvectors. Each eigenvector can have up to  $k$  non-zero elements, but note that in practice the sparseness of the transition probability matrix leads to the eigenvectors being sparse as well.

**[0070]** The cost of computing the similarity of a  $k$  component profile against an  $n$  word target text is as follows. It takes  $O(n)$  operations ( $O$  denotes big- $O$  the order of magnitude

of the operation) to tokenize and parse the target text. Each operation is relatively inexpensive, and consists of forming new tuples and tokens as they are encountered. In addition, tally the occurrence of each tuple-token pair. The result of parsing is a sparse representation of the transition probability matrix.

**[0071]** Computing the similarity entails  $2 \cdot k$  vector-matrix operations, where each vector matrix operation is a  $k$ -vector against a  $k$ -by- $n$  sparse matrix, or  $O(2 \cdot k \cdot k)$  operations, or simply  $O(k^2)$ . Notice that because the  $k$ -by- $n$  matrix is sparse, a non-zero term only arises when there is a tuple in the eigenvector that matches a tuple in the matrix, or when there is a token in the eigenvector that matches a token in the matrix.

**[0072]** Adding the contribution of the parsing cost  $O(n)$  and the similarity computation cost  $O(k^2)$ , one gets the overall cost of  $O(n) + O(k^2)$  to compute a single similarity measure.

**[0073]** In practice, it is possible to reuse the results of the parsing step for several similarity computations. For example, parse an email, contract, or patent document once to compute the sparse transition matrix, and apply the similarity computation for many candidate profiles to determine the best match. In other words, the cost to run an  $n$ -token target text against  $m$  possible profiles, each of which uses  $k$  components has the following cost:

$$\text{Cost of computing similarity} = O(n) + O(m \cdot k^2),$$

Where  $n$  = number of tokens in the target input

$m$  = number of profiles to compare against



k = number of components in each profile.

## **ADDITIONAL APPLICATIONS**

### **Term value extraction**

**[0074]** Discussed below is an example of one embodiment of the present invention for term value extraction.

**[0075]** Even when contracts are executed in electronic form or are transformed into blocks of machine-readable text using OCR techniques, there remains the problem of extracting the values of terms contained within it. What is the effective date of the contract? When does the contract expire? Who is the supplier? A trained human would have little difficulty answering these questions. Multiply the problem by the hundreds of suppliers that typically serve a large organization and the myriad of existing and expired contracts that reside on file. Consider also the problem of two companies merging their operations, a straightforward business query asks to identify all the common suppliers with expiration dates that have not yet come due. A manual solution to this problem entails significant time and expense, essentially becoming a huge data entry undertaking. Automatic extraction of term values within the raw contract text saves time and money.

**[0076]** In one embodiment of the present invention, this example of an algorithm may be used to extract term values from a block of text.

1. Begin with one or more examples where the term value appears. Use these as the training set.

2. Use the profiling algorithm described earlier in a special way as described here.

Choose a window of  $w$  words that precedes the term value. Consider the following marker for the effective date in the introductory clause in a contract: "[effective-date]". Replace the marker, "[effective-date]", with the token TT\_TV. Tokenize  $w$  words that precede the marker, and  $w$  words that follow this marker. For example, if  $w=4$ , then tokenize the following text block:

"is effective on the TT\_TV ( effective date )"

The token TT\_TV is used to denote the term value itself. Similarly, "(" and ")" are used as tokens to represent the open and close parenthesis, respectively.

3. Compute and save the profile for the text block described above.
4. To extract the effective date from a block of text, perform the following search over the text block, trying out different candidate term values. In particular, a candidate term value can be thought of as starting at the  $s$ -th starting token, and proceeding for  $t$  tokens. For example, the term value candidate (5,2) consists of the two tokens commencing on the 5<sup>th</sup> token - in other words, the 5<sup>th</sup> and 6<sup>th</sup> tokens.
5. For each possible candidate (s,t) pairs, where  $s$  varies from 1, corresponding to the first token, to the  $N$ th token in the target block, and  $t$  varying from a length of 1 up to the maximum length the term value is assumed to be, compute the similarity score for each such (s,t).

6. Choose the (s,t)\* that gives the largest similarity score.
7. Extract the term value by starting from the s-th token and selecting the next t tokens.

### Example

**[0077]** Here is an example of some training clauses (a, b, and c).

- (a) This Materials and/or Equipment Contract is effective on the [effective-date] ([Effective Date] [buyer] a corporation organized and existing under the laws of the state of [buyer-state] (hereinafter [Company] and [seller] a corporation organized and existing under the laws of the state of [seller-state] (hereinafter [Supplier].. [source: Master Materials and/or Equipment]
- (b) This Purchase Agreement shall be effective [effective-date]; is made ([Agreement Date], by and between [buyer]. ([Buyer], a [buyer-state] corporation, having principal offices at [buyer-address] and [supplier], having principal offices at [supplier-address].. [source: AGREEMENT Purchasing Short]
- (c) This Volume Purchase Agreement ([VPA] is made this [effective-date] by and between [buyer], having principal offices at [buyer-address] ([Buyer] and [supplier] ([Supplier], having principal offices at [seller-address].. [source: volume purchase agreement]

**[0078]** Shown in Figure 5 is a target clause with the desired business term values

underlined.

**[0079]** Shown in Table 17 is what the algorithm obtained as the extracted business term values.

Table 17. Results of the term extraction procedure.

Effective date	22nd day of January 2003 □ □
Seller	Advanced Tire Products Inc a □

### **QUICK SIMILARITY SEARCHING**

**[0080]** Discussed below are several embodiments of the invention illustrating similarity searching. The present invention is not limited to only these embodiments.

**[0081]** Assume that a contract has been transformed into machine-readable form. Each clause now needs to be matched with its closest counterpart in an existing library of contract clauses. In one embodiment of the invention, a straightforward, albeit brute-force, approach to this problem may be used which computes the similarity score of each of the M clauses in the target contract against each of the N clauses in a standard clause library.

**[0082]** In another embodiment of the invention the approach here may be used to solve this search problem in substantially  $M \cdot \log(N)$  time, instead of the  $M \cdot N$  time of the straightforward approach. This approach involves pre-computing an indexing structure around the profiles and proceeds thusly:

1. Start with N text blocks.

2. Given a target text block T. The goal is to find the closest text block of T among the N text blocks.
3. Concatenate the N blocks to form a single text block.
4. Compute the profile of the single block.
5. Obtain the  $C \leq N$  eigenvalues from the profile, and sort them by their magnitude. Call them  $w_0, \dots, w_{C-1}$ , where  $w_0$  is the eigenvalue with the largest magnitude.
6. For each of the N text blocks  $t_i$ , compute the transition probability matrix,  $b_i$ .
7. Starting from the largest eigenvalue  $w_0$ , compute the partial sum,  $s_i$  for each text block.

$$s_i = u_{\text{transpose}_0} * b_i * v_0$$

8. Each of the  $s_i$  is a complex number in the complex plane. Partition the complex plane down the imaginary axis, the half with non-negative real part and the other half with negative real part. Each of the  $s_i$  will fall on one side or the other. Count the fraction  $f$  that falls on the non-negative real part side. One may say that there is a sufficiently good partitioning of the  $n$  points, if the fraction  $f$  falls within a predetermined range  $[0.3, 0.7]$  for example. If there is a good partitioning, then select left and right eigenvalue pair corresponding to  $w_0$  as the partitioning component.

9. If  $w_0$  does not yield a good partitioning, then proceed to the eigenvector with the next largest magnitude. Do this iteratively until the  $k$ -th left/right eigenvector pair that yields a good partition is found. This is the partitioning component.
10. In practice, there will always be a largest component  $k$ . This follows from the properties of how the principal component analysis technique chooses the eigenvectors in the first place, or this follows from the definition of the left and right eigenvectors themselves.
11. Once the partitioning component is determined, use it to partition the  $n$  text blocks into the each respective half-plane.
12. Take each subset of text blocks and repeat the steps **Error! Reference source not found.-Error! Reference source not found.** above.
13. Repeat this procedure on each subset until each subset contains only one text block. In essence the final partition computes the specific profile for that sole text block  $T$ .

**[0083]** In another embodiment of the present invention, the technique used for searching relies on the observation that computing a profile is relatively more expensive than computing similarity. Therefore, an approach that computes a profile once and uses it many times is introduced. Computing similarity of a given profile over a larger block of text entails computing the probability transition matrix over the text. In practice the probability transition matrix is sparse, which keeps the memory requirement from growing too rapidly as the text block size increases.

This algorithm exploits the ability of a profile to be used repeatedly over a collection and so is referred to as a "shuffle" search.

1. Start with N text blocks.
2. Assume a given target text block T. The goal is to find the block among the N that is most similar to T. Compute a profile for T.
3. Partition the N text blocks into equal sized sets. There are known "perfect shuffle" algorithms that may accomplish this. Concatenate the text of a given partition into a single block. For convenience, call the resulting blocks red and black. Run the profile for T over red and black. Determine the set that has the higher similarity. Tally a "vote" for each member of the set scoring the higher similarity.
4. Pick a different partition of the N text blocks. Repeat step 0 until a predetermined number of votes are taken. A good approach is to take more votes for larger size N, such as  $\log(k*N)$ , where k is a small fudge factor.
5. When the votes have completed, sort the N text blocks according to the number of votes that each obtained. Pick the top fraction f of N, where  $0 < f < 1$ . A value like  $f=0.5$  or  $f=0.75$  works well. Eliminate the bottom  $1-f$  portion of the N text blocks.
6. Repeat steps 0-0, until all but a small number of K text blocks have been eliminated, where K is a small number like 4 or 8.

7. For each of the K remaining text blocks compute the similarity using the profile of T.
8. Sort by similarity. The result is the closest matching text blocks to the target block T.

This technique searches a set of size N in approximately  $\log(N)$  time. As noted previously, it has the additional advantage that the profile of the target T is computed once and used many times.

**[0084]** This technique has another advantage when the same set of N blocks is searched repeatedly with different targets, which is common in practice. For example, the N blocks might be a list of N company names and incoming blocks of text are different fragments of company names, possibly with misspellings, obtained via optical character recognition. One optimization is to cache the partitions and sub-partitions of the N. In other words, as the search targets T are processed, different partitions of N will be computed. Computing the similarity entails constructing the terms of the probability transition matrix for the concatenated text blocks corresponding to a given partition. Instead of discarding the transition matrix values, save them according to the partition. When memory is finite, a good scheme may be to define a memory quota and to save the most recently used transition matrix values for their partitions. This approach will minimize the number of times that the probability transition matrix will need to be recomputed from scratch.

**[0085]** As an example, the shuffle search technique was used to obtain the most similar



matching pages in Table 16, shown earlier.

## TEXT SIMILARITY

**[0086]** In one embodiment of the present invention it may be used as described above to identify similar blocks of text. For example, in the operation of a company it is common to have legal contracts with customers and suppliers. After a contract has been executed, it is common for the paper contract to sit untouched in a filing cabinet. But during the lifetime of the contract, business situations change which may necessitate reviewing the terms of the contract. A company may find it necessary to audit its existing contracts to determine which of them use a particular kind of liability clause, or a certain wording of a termination clause, say, for its suppliers. Or, a company merges with another, and there is suddenly an overlap of contracts between the merging companies with their common suppliers. Finding a particular clause is difficult when the specific wording varies. For example, the word, "termination" appears in multiple places in a contract. Of those occurrences, a smaller handful of clauses contain the relevant termination clause. Different section titles mark the clause: "Termination of an order" or "Termination of request for services." Even within the body of the clause, one clause may use the phrase, "at any time, upon issuance of written notice," while another may state, "at any time, without cause, terminate an Order in whole or in part upon written notice."

In one embodiment of the invention, this procedure may be used for identifying the paragraphs in a document, viewed as a collection of text blocks:

1. Assemble paragraphs as text blocks to form the paragraph library.
2. For each text block in the library, compute a profile and save it.
3. Construct a similarity tree as defined previously.
4. Given a target contract to identify, break it into clauses by manual or automatic means. Manual means may be that a person selects the points where one paragraph (clause) ends and the next start. Automatic means may be that the paragraph and section formatting information is used to deduce the end of one clause and the beginning of the next clause.
5. For each paragraph (clause) in the target document, use one of the searching approaches described above to identify the most closely matching clauses.

## **GENERALIZATIONS AND EXTENSIONS**

**[0087]** One of skill in the art will appreciate that the prior examples and prior discussion were limited in scope and applicability in order to convey the essence of the present invention in a form that was readily understood. The present invention is not so limited, and so here generalization and extensions of the present invention are discussed.

**[0088]** Any observable process whose behavior can be represented by a sequence of tokens can be parsed and recognized using the techniques described here. It can be seen that a block of text is a special case of a sequence of tokens.

## Applications of technique

### HTML/XML as a sequence of tokens

**[0089]** For example, a common problem occurs in migrating web content that entails extracting content-rich portions of HTML into XML. It is typical for blocks of content to be represented as tables, or tables nested within tables. Web pages often use tables to create sidebars or related links. The challenge that arises in extracting this content is to recognize occurrences of such tables. An HTML table is a rooted tree of elements. One can apply the similarity techniques described earlier if an HTML table can be converted into a sequence of tokens. Here is one way to do this. Tokenize the HTML beginning from the root. Tokenize each HTML element as a special token—for example, <table> is parsed as the token "H-table" and <tr> becomes "H-tr." Tokenize the text contained within and between elements in the same way as normal text parsing. For example, use the techniques described earlier. All closing HTML elements such as </td> and </table> are tokenized as H-END. Consider an example.

1	2
3	4

Example of an HTML table.

In HTML, this is written as,

```
<table>
  <tr><td>1</td><td>2</td></tr>
  <tr><td>3</td><td>4</td></tr>
</table>
```

After tokenization, this is

H-table

```
      H-tr H-td 1 H-END H-td 2 H-END H-END
      H-tr H-td 3 H-END H-td 4 H-END H-END
H-END
```

Notice that the whitespace does not matter. To make it easy to compare with the original HTML, the same indentation structure was used. Notice how the common H-END token denotes the closing of an HTML element.

By using this approach, HTML similarity can be computed in exactly the same way that the earlier example computed text similarity.

### **Observed behavior as sequence of tokens**

**[0090]** To illustrate the point that many kinds of problems can be recast as a sequence of tokens, if one wanted to study how someone flips channels while watching television, one can represent their behavior as, "What channel am I tuned to at time t?" One can interpret turning the television off as completing a sentence. A new sentence begins when the television is turned on.

**[0091]** In this application there would be profiles to describe different modes of television watching. There is the profile associated with, "I'm really captivated by this movie and don't change the channel at all." Another profile might be, "I'm interested, but every time there's a commercial break, I surf off to another channel to see what's going on." Yet another profile might be, "I'm really bored, and cannot find anything to watch, so I continuously surf through all the movies that are currently playing."

## Extensions

**[0092]** In practice, the Applicant has found that the procedure described above has excellent performance, correctly representing textual similarity in numerical terms. Upon deeper examination of the procedure, it will be noted that it is possible to expand the way that the probability transition matrix  $a[]$  is computed. For instance, when a history window of  $h$  tokens is used to compose a tuple, and the probabilistic transition to the subsequent token, a model of  $(h+1)$ -gram token sequences is built. Given a current tuple, the model accounts for the possible transitions to the next word. One possible disadvantage of such a formulation may be the precision itself. The phrases, "read in pickle color," and "read in color" are entirely different, for example. One approach to compensate for this over-preciseness may be to widen the transition window. Instead of tallying the words that fall exactly in the next position following a given tuple, a somewhat more "forgiving" model may tally all the words that fall within  $r$  tokens of the originating tuple. When  $r=1$  then the model described earlier results. For  $r>1$  one can capture the idea that the transition from the tuple "read in" to the token "color" is essentially equivalent in the phrases, "read in color" and "read in pickle color."

**[0093]** By introducing the parameter  $r$  for the size of the window to tally transitions from tuple to token, a continuum of models for the underlying token sequence is introduced. As already noted  $r=1$  yields the original model described. Going to the other extreme of very large  $r$ , one will note that in the limit  $r \rightarrow \infty$  encompasses the entire text. And if the text has essentially the same word frequencies throughout, then observe that in the limit, all rows of the transition matrix will tend to converge to the frequency distribution of the words that

comprise the text. In other words, the classic word frequency model has been derived.

One of skill in the art will appreciate how novel it is to define a general model as was done here that embeds the stationary Markov n-gram model on one extreme ( $r=1$ ), and the word frequency model at the other extreme ( $r=\infty$ ). Between these extremes, the present invention in various embodiments can adjust the model to represent locality at the phrase level ( $r=4$ , for example), or at the sentence level ( $r=20$ , for example), or at the paragraph level ( $r=200$ , for example). Small values of  $r$  will produce very discriminating similarity profiles, while large values of  $r$  will yield broad-brush similarity measures.

**[0094]** Nor is this the full range of embodiments of the present invention.

**[0095]** Other embodiments of the present invention are possible. For example, notice that a token transition window of width  $r$ , as discussed earlier, is substantially a step function of width  $r$  that assigns an equal weight  $1/r$  to a token at the 1, 2,  $\dots$ ,  $r$  locations. In other embodiments of the invention, it is possible to give greater weight to nearby tokens compared to farther away tokens, by using a transition weight function  $s(i)$ , where  $0 \leq s(i) \leq 1$ , for  $i=1, \dots, r$ , and normalized so that  $\sum_{i=1}^r s(i) = 1$ .

**[0096]** What is to be appreciated is that embodiments of the present invention are applicable to ANY observable process whose behavior can be represented by a sequence of tokens.

**[0097]** Figure 6 illustrates one embodiment of the invention in flow chart form. Two columns of flow are indicated, one for Reference, the other for Target. Reference input is received at 602, at 604 the reference input received at 602 is tokenized. At 606 a transition probability model of the tokens from 604 is generated. At 608 an eigenspace is generated

based on the transition probability model from 606. Target input is received at 612, at 614 the target input received at 612 is tokenized. At 616 a transition probability model of the tokens from 614 is generated. At 620 input is received from the eigenspace at 608 and from the transition probability model at 616 and a similarity measure is generated.

**[0098]** In other embodiments, a profile at 608 as represented by an eigenspace may be pre-computed and used to generate a similarity measure 620 against a Target input. In other embodiments, as previously described a shuffle operation may occur, for example, between blocks 614 and 616.

**[0099]** Figure 7 through Figure 10 illustrate embodiments of the invention showing two different possible search approaches (denoted Search I and Search II). Both Search I and Search II address a similar problem. Namely, assume that there are NO reference text blocks that are known in advance. If one is then given a target text block T, the problem is to find the most similar block from the N in  $O(\log(N))$  time.

**[00100]** Search I as illustrated in Figure 7, Figure 8, and Figure 9 entails two main phases. In the first phase, is built a binary tree which indexes the text blocks. In the second phase, knowing the target T, one traverses the tree to find the closest matching block.

**[00101]** Figure 7 illustrates building a binary tree. The idea is to use an left/right eigenvector pair within a profile to partition the set of N blocks into  $f \cdot N$  blocks, where f is in the interval [0.3, 0.7]. In general, one wants it to be in  $0.5 \pm \delta$ , for some delta. So the approach is to concatenate all N text blocks, and to compute a profile of that. Obtain the eigenvalues and arrange them in decreasing magnitude. Starting from the eigenvalue with

the largest magnitude, obtain the corresponding left and right eigenvectors  $u$  and  $v$ . Iterate through the  $N$  blocks and compute the transition matrix  $B$ . Compute  $u \cdot \text{hermitian} * B * v$ , which gives a complex number. Use the real component of this complex number to consider a particular block as in the right-half (real part  $\geq 0$ ), or as in the left-half (real part  $< 0$ ). Compute  $f$ . If  $f$  is in the interval  $0.5 \pm \delta$ , then use the  $u$  and  $v$  as the separating eigenvectors. Otherwise, go on to the next smaller eigenvalue. Eventually, a suitable  $u$  and  $v$  are found. This partitions the set of  $N$  blocks into two smaller sets.

**[00102]** The approach above is applied to each of the left and right sets. Eventually, the result is singleton sets. In another embodiment, one may say that for some threshold, say 4, to just individually compute the similarity, instead of using the partitioning approach. In either approach the net result is that a binary tree has been built.

**[00103]** Figure 8 illustrates one embodiment of the invention to obtain partitioning eigenvectors. This approach was described above.

**[00104]** Figure 9 illustrates one embodiment of the invention showing how to use a binary tree. Given a binary tree built as described earlier, the problem is to find the closest matching reference block. Start from the root, and apply the partitioning eigenvectors at each node in the tree to the transition matrix of the target block  $T$ . Use that to the real part of the product  $u \cdot \text{hermitian} * B \cdot T * v$  to decide whether to go left or right. Continue recursively, until a node is found that just holds one text block. That is the closest match.

**[00105]** Figure 10 illustrates one embodiment of the invention showing a search approach denoted Search II. This approach uses the same idea as in Search I, except that instead of computing profiles for the reference blocks, one instead computes a profile of a



target block T. For the profiles, use the transition matrices of different aggregations of the N blocks. Take N blocks, shuffle them into two sets; call them A and B. Concatenate each set. Compute a transition probability matrix for each set. Compare T's profile against A and B. If A has a higher similarity score, tally a vote for each member of A. Same thing if B has a higher similarity score. Repeat the random shuffles some number of times R. R is typically a small number like 4 or 8. Or R can be determined dynamically, by seeing when clear winners begin to emerge.

**[00106]** Once the voting has completed, retain the highest vote getters. Say the N/2 reference blocks receive the most votes. Now repeat the procedure using the smaller set. Eventually, one gets K surviving blocks. These are the K most similar blocks. In another embodiment, one can compute the similarity of each of the K to get a detailed ranking of the K blocks.

**[00107]** Figure 11 through Figure 35 illustrate one embodiment of the invention as a demonstration. Figure 11 is the title page for demo walkthrough. Figure 12 states one usage of text similarity is to help import contracts into a management system. Figure 13 illustrates the flow of contract content, from paper and electronic softcopy into a contract management repository. The repository becomes a source of controlled, indexed clauses which are emitted into a contract document. After editing, the contract can be reimported into the repository. Figure 14 shows a sample contract that might be imported into the contract management system. Figure 15 summarizes some of the problems with current approaches. Figure 16 shows how conventional comparison tools do not handle differences well, such as formatting shown here. Figure 17 shows that using a comparison

tool after pre-converting both inputs to plain unformatted text doesn't improve the comparison.

**[00108]** Figure 18 illustrates how one embodiment of the invention may provide a solution that uses a library of contracts both as whole contracts and as individual clauses. Profiles are computed based on the whole contracts and on the clauses. Figure 19 shows a sample collection of contracts in the library of contracts, which forms the training set. Figure 20 shows how to generate profiles based on the contracts in the contract library. Figure 21 shows how to break contracts down into individual clauses, and construct a library of clauses. Figure 22 shows how to generate profiles based on the clauses in the clause library. Figure 23 shows how to "recognize" a target contract, which is generally a contract not in the contract library, and compute the similarity of the target against pre-computed contract profiles.

**[00109]** Figure 24 illustrates an overview that summarizes the main steps of the procedure: compute similarity of the target at the contract level, and then compute similarity of individual clauses of the target at the clause level. Figure 25 illustrates a target contract being processed as a whole and as individual clauses. Figure 26 has a graph showing the numerical similarity score of contracts in the library against the target. Figure 27 illustrates the target contract and its closest match in the library. Figure 28 shows that the matching contract in the library has a different internal structure from the target contract. Figure 29 has a graph showing the numerical similarity scores of clauses from the clause library relative to one of the target clauses. Figure 30 shows the matching clauses - note that neither is a word-for-word match for the other. Figure 31 shows another clause from the

clause library that had a good match against the target clause.

**[00110]** Figure 32 illustrates a different clause from the target contract, and the numerical similarity scores of clauses from a contract in the contract library. Figure 33 has a detailed view of the most similar clause in the clause library, where 1, 2, 3, and 4 are labels for the clauses. Figure 34 shows a comparison of the previous figure with the target clause where 1, 2, 3, and 4 are labels for clauses that most correspond to those of the previous figure. Figure 35 summarizes the main points - build a contract library and compute similarity of a target at the contract level and at the clause level.

**[00111]** Figure 36 is a cover sheet from a scenario walkthrough. Figure 37 shows a clause that has been annotated with markers for the locations where business term values (BTV) occur as does Figure 38. Figure 39 shows the overall flow of data-annotated training documents on the left, input contracts on the right, which are processed by the recognition engine to compute similar clauses and to extract BTV's.

**[00112]** Figure 40 is a cover sheet from a overview of the business term value (BTV) extraction technology. Figure 41 indicates the internal status of various embodiments of the invention as of 1/24/2002. Figure 42 illustrates the overall process flow of the BTV procedure. Figure 43 shows an example of a clause that has been annotated with BTV markers (also referred to as "tags"). Figure 44 shows a training set for BTVs that occur within the introduction clause in a contract. Figure 45 is an example of a target clause - the goal will be to extract the BTVs contained within this clause. Figure 46 shows the results of the BTV extraction for the "effective date." Figure 47 is another example of a target clause - the goal is to extract the supplier name. Figure 48 shows the annotated clauses for the

"agreement term" BTV. This is the training set. Figure 49 shows the results of the BTV extraction for the "agreement term." Figure 50 lists further embodiments of the invention as of 01-24-03.

**[00113]** Figure 51 is a cover sheet for a contract analysis walkthrough. Figure 52 shows the results of comparing two contracts at the clause level, showing the pairs of clauses that have high similarity scores. Figure 53 is a graphical representation of the matching pairs of clauses. Figure 54 is a rank order plot of the clause similarities. Figure 55 shows a clause from the services contract - it is the warranty clause. Figure 56 shows the most similar clause from the materials contract, relative to the previously shown clause. Figure 57 shows another clause from the materials contract that was also similar to the warranty clause in the services contract.

**[00114]** Figure 58 shows another clause from the services contract - it contains the term and termination clauses. Figure 59 shows the most similar clause from the materials contract. Figure 60 shows a clause from the services contract which is similar. Figure 61 shows a similar clause from the materials contract.

**[00115]** Figure 62 is an example of a target clause for which a BTV extraction will be done. Figure 63 shows the results of the BTV extraction: effective date and seller. Figure 64 is an example of a training clause for the effective date, buyer, seller, and seller state BTVs. Figure 65 shows more examples of training clauses.

**[00116]** Figure 66 shows an overview of a procedure to compute similarity of two text blocks. Figure 67 illustrates one way to tokenize. Figure 68 shows an example of an input text tokenized. Figure 69 shows a procedure to compute a profile. Figure 70 shows the

tokens and tuples that result from the n-gram and probability state transition analysis.

Figure 71 shows other operations that may be performed.

**[00117]** Figure 72, Figure 73, and Figure 74 is an XML representation of a profile in one embodiment of the invention. It has sections showing the definitions of the tokens, the tuples, the eigenvalues, the left eigenvectors, and the right eigenvectors. At 1 is a token definition section with 6 tokens (v id="0" through "6"). For example, token "1" is TT\_TERM. At 2 is a tuples definition section with 7 tuples (V id="0" through "6"). For example, tuple "4" is show.our. At 3 is an indicator that there are 7 components (10, and others starting at 20, 30 (Figure 73), 40 (Figure 73), 50 (Figure 74), 60 (Figure 74) and 70 (Figure 74)). The first component 10 has an eigenvalue 7, left eigenvectors 5, and a right eigenvectors 6 associated with it.

**[00118]** At 7 is an eigenvalue and an indication that there are 9 parts to this first component 10 (component count="9"). Eigenvalue 7 has an imaginary value (eim="0.0"), a magnitude value (ema="1.0000000000000018"), and real value (ere="1.0000000000000018").

**[00119]** Left eigenvectors 5 has five eigenvectors (v id="0", "5", "3", "4", and "6"). Each eigenvector has several parameters listed. For example, at 8 a single eigenvector has:

v id="0"	A tuple identification
im="0.0"	The imaginary part of the eigenvector
ma="0.4472135954999585"	The magnitude of the eigenvector
ph="-3.141592653589793"	The phase of the eigenvector (in radians)

re="-0.4472135954999585"	The real part of the eigenvector
sym="TT_TERM.TT_TERM."	What the tuple identification represents

**[00120]** Right eigenvectors 6 has four eigenvectors (v id="0", "4", "3", and "5"). Each eigenvector has several parameters listed similar to the left eigenvectors.

**[00121]** Figure 75 is an XML representation of a null profile.

**[00122]** Figure 76, Figure 77, Figure 78, Figure 79, and Figure 80 is another example of an XML representation of a profile.

**[00123]** Thus a method and apparatus for fundamental operations on token sequences: computing similarity, extracting term values, and searching efficiently have been described.

**[00124]** Figure 1 illustrates a network environment 100 in which the techniques described may be applied. The network environment 100 has a network 102 that connects S servers 104-1 through 104-S, and C clients 108-1 through 108-C. More details are described below.

**[00125]** Figure 2 illustrates a computer system 200 in block diagram form, which may be representative of any of the clients and/or servers shown in Figure 1, as well as, devices, clients, and servers in other Figures. More details are described below.

**[00126]** Referring back to Figure 1, Figure 1 illustrates a network environment 100 in which the techniques described may be applied. The network environment 100 has a network 102 that connects S servers 104-1 through 104-S, and C clients 108-1 through 108-C. As shown, several computer systems in the form of S servers 104-1 through 104-S and C clients 108-1 through 108-C are connected to each other via a network 102, which

may be, for example, a corporate based network. Note that alternatively the network 102 might be or include one or more of: the Internet, a Local Area Network (LAN), Wide Area Network (WAN), satellite link, fiber network, cable network, or a combination of these and/or others. The servers may represent, for example, disk storage systems alone or storage and computing resources. Likewise, the clients may have computing, storage, and viewing capabilities. The method and apparatus described herein may be applied to essentially any type of communicating means or device whether local or remote, such as a LAN, a WAN, a system bus, etc.

**[00127]** Referring back to Figure 2, Figure 2 illustrates a computer system 200 in block diagram form, which may be representative of any of the clients and/or servers shown in Figure 1. The block diagram is a high level conceptual representation and may be implemented in a variety of ways and by various architectures. Bus system 202 interconnects a Central Processing Unit (CPU) 204, Read Only Memory (ROM) 206, Random Access Memory (RAM) 208, storage 210, display 220, audio, 222, keyboard 224, pointer 226, miscellaneous input/output (I/O) devices 228, and communications 230. The bus system 202 may be for example, one or more of such buses as a system bus, Peripheral Component Interconnect (PCI), Advanced Graphics Port (AGP), Small Computer System Interface (SCSI), Institute of Electrical and Electronics Engineers (IEEE) standard number 1394 (FireWire), Universal Serial Bus (USB), etc. The CPU 204 may be a single, multiple, or even a distributed computing resource. Storage 210, may be Compact Disc (CD), Digital Versatile Disk (DVD), hard disks (HD), optical disks, tape, flash, memory sticks, video recorders, etc. Display 220 might be, for example, a Cathode Ray

Tube (CRT), Liquid Crystal Display (LCD), a projection system, Television (TV), etc. Note that depending upon the actual implementation of a computer system, the computer system may include some, all, more, or a rearrangement of components in the block diagram. For example, a thin client might consist of a wireless hand held device that lacks, for example, a traditional keyboard. Thus, many variations on the system of Figure 2 are possible.

**[00128]** For purposes of discussing and understanding the invention, it is to be understood that various terms are used by those knowledgeable in the art to describe techniques and approaches. Furthermore, in the description, for purposes of explanation, numerous specific details are set forth in order to provide a thorough understanding of the present invention. It will be evident, however, to one of ordinary skill in the art that the present invention may be practiced without these specific details. In some instances, well-known structures and devices are shown in block diagram form, rather than in detail, in order to avoid obscuring the present invention. These embodiments are described in sufficient detail to enable those of ordinary skill in the art to practice the invention, and it is to be understood that other embodiments may be utilized and that logical, mechanical, electrical, and other changes may be made without departing from the scope of the present invention.

**[00129]** Some portions of the description may be presented in terms of algorithms and symbolic representations of operations on, for example, data bits within a computer memory. These algorithmic descriptions and representations are the means used by those of ordinary skill in the data processing arts to most effectively convey the substance of their work to others of ordinary skill in the art. An algorithm is here, and generally, conceived to



be a self-consistent sequence of acts leading to a desired result. The acts are those requiring physical manipulations of physical quantities. Usually, though not necessarily, these quantities take the form of electrical or magnetic signals capable of being stored, transferred, combined, compared, and otherwise manipulated. It has proven convenient at times, principally for reasons of common usage, to refer to these signals as bits, values, elements, symbols, characters, terms, numbers, or the like.

**[00130]** It should be borne in mind, however, that all of these and similar terms are to be associated with the appropriate physical quantities and are merely convenient labels applied to these quantities. Unless specifically stated otherwise as apparent from the discussion, it is appreciated that throughout the description, discussions utilizing terms such as "processing" or "computing" or "calculating" or "determining" or "displaying" or the like, can refer to the action and processes of a computer system, or similar electronic computing device, that manipulates and transforms data represented as physical (electronic) quantities within the computer system's registers and memories into other data similarly represented as physical quantities within the computer system memories or registers or other such information storage, transmission, or display devices.

**[00131]** An apparatus for performing the operations herein can implement the present invention. This apparatus may be specially constructed for the required purposes, or it may comprise a general-purpose computer, selectively activated or reconfigured by a computer program stored in the computer. Such a computer program may be stored in a computer readable storage medium, such as, but not limited to, any type of disk including floppy disks, hard disks, optical disks, compact disk- read only memories (CD-ROMs), and

magnetic-optical disks, read-only memories (ROMs), random access memories (RAMs), electrically programmable read-only memories (EPROM)s, electrically erasable programmable read-only memories (EEPROMs), FLASH memories, magnetic or optical cards, etc., or any type of media suitable for storing electronic instructions either local to the computer or remote to the computer.

**[00132]** The algorithms and displays presented herein are not inherently related to any particular computer or other apparatus. Various general-purpose systems may be used with programs in accordance with the teachings herein, or it may prove convenient to construct more specialized apparatus to perform the required method. For example, any of the methods according to the present invention can be implemented in hard-wired circuitry, by programming a general-purpose processor, or by any combination of hardware and software. One of ordinary skill in the art will immediately appreciate that the invention can be practiced with computer system configurations other than those described, including hand-held devices, multiprocessor systems, microprocessor-based or programmable consumer electronics, digital signal processing (DSP) devices, set top boxes, network PCs, minicomputers, mainframe computers, and the like. The invention can also be practiced in distributed computing environments where tasks are performed by remote processing devices that are linked through a communications network.

**[00133]** The methods of the invention may be implemented using computer software. If written in a programming language conforming to a recognized standard, sequences of instructions designed to implement the methods can be compiled for execution on a variety of hardware platforms and for interface to a variety of operating systems. In addition, the

present invention is not described with reference to any particular programming language. It will be appreciated that a variety of programming languages may be used to implement the teachings of the invention as described herein. Furthermore, it is common in the art to speak of software, in one form or another (e.g., program, procedure, application, driver,...), as taking an action or causing a result. Such expressions are merely a shorthand way of saying that execution of the software by a computer causes the processor of the computer to perform an action or produce a result.

**[00134]** It is to be understood that various terms and techniques are used by those knowledgeable in the art to describe communications, protocols, applications, implementations, mechanisms, etc. One such technique is the description of an implementation of a technique in terms of an algorithm or mathematical expression. That is, while the technique may be, for example, implemented as executing code on a computer, the expression of that technique may be more aptly and succinctly conveyed and communicated as a formula, algorithm, or mathematical expression. Thus, one of ordinary skill in the art would recognize a block denoting  $A+B=C$  as an additive function whose implementation in hardware and/or software would take two inputs (A and B) and produce a summation output (C). Thus, the use of formula, algorithm, or mathematical expression as descriptions is to be understood as having a physical embodiment in at least hardware and/or software (such as a computer system in which the techniques of the present invention may be practiced as well as implemented as an embodiment).

**[00135]** A machine-readable medium is understood to include any mechanism for storing or transmitting information in a form readable by a machine (e.g., a computer). For

example, a machine-readable medium includes read only memory (ROM); random access memory (RAM); magnetic disk storage media; optical storage media; flash memory devices; electrical, optical, acoustical or other form of propagated signals (e.g., carrier waves, infrared signals, digital signals, etc.); etc.

**[00136]** As used in this description reference to a transition probability matrix is referring to a transition probability model which is represented in the form of a matrix. To succinctly denote this, transition probability matrix is often used instead of calling it a transition probability model represented as a matrix. Thus, the term transition probability matrix is used to denote a transition probability model which may be represented conveniently as a matrix, however it is not limited to only a matrix representation. What is to be understood is that as used with tokens, the transition probability matrix is a model of a token to another token transition sequence.

**[00137]** As used in this description reference to "eigenspace," "Eigenspace," or similar phrases with reference to a profile of  $k$  components refers to  $k$  eigenvalues,  $k$  left eigenvectors, and  $k$  right eigenvectors for such a profile. Thus, the term eigenspace is used to denote the various eigenvalues and eigenvectors associated with a profile. What is to be appreciated is that an eigensystem analysis of a transition probability model transforms the state space of tokens into the eigenspace. As one of skill in the art knows such a transform has admirable qualities such as orthogonal basis elements.

**[00138]** As used in this description reference to "eigensystem analysis," "solving the eigenproblem," "singular value decomposition," or similar phrases are to be understood as performing those operations on an eigenspace or space represented in another form so as

to yield the results discussed. The terms are used interchangeably and their meaning will be evident to one of skill in the art based on the surrounding context. For example, one of skill in the art will appreciate that solving a square matrix can be formulated as an eigenproblem, while solving a non-square matrix can be formulated as a singular value decomposition problem. In embodiments of the present invention, whether the transition matrix is square or non-square depends on whether the number of tokens happens to equal the number of tuples, and by how much one may reduce dimensionality for practical considerations.

**[00139]** As used in this description reference to "component", "C", or similar phrases with respect to eigenvectors pairs refers to entities that may make up the component. These entities are a real eigenvalue and the corresponding left and right eigenvector pair, or a complex conjugate pair of eigenvalues and a corresponding pair of complex left and right eigenvectors. Thus, an  $i$ th component may be represented by  $(w.i, u.i, v.i)$ , where  $w.i$  is the  $i$ th eigenvalue(s),  $u.i$  is the left eigenvector(s), and  $v.i$  is the  $i$ th right eigenvector(s).

**[00140]** What is to be appreciated is that depending upon the granularity or level of recognition or similarity that one wants, the present invention by allowing a range of arbitrarily tokenized inputs allows for similarity measurements. For example, if an entire document is to be classified as to a major category, for example, Non-Disclosure Agreement (NDA), Purchase Order (PO), or Contract (K) then the target document (X) would be compared to the NDA, PO, and K profiles by the methods described. If the classification is to be, for example, by paragraph, then paragraph profiles of the NDA, PO, and K would be compared to the target paragraphs. What is to be appreciated is that the

profiles may be pre-computed and used many times against the more easily generated transition probability model of the target. Additionally, the paragraphs or clauses of the profiled documents and the targets need not be the same nor does the tokenization process or the windowing for the profiles or references need to be the same as for the targets. For convenience in describing the invention, the phrase "clause" may be used to indicate an entire document, a paragraph, or any group of information. Additionally, for ease in understanding the invention, the phrase "text" while often used to indicate actual alphanumeric strings, is not so limited and may more generally be considered any data.

**[00141]** The description of the invention in the discussion has illustrated how to locate a high degree of similarity between a profile and a target, one of skill in the art will recognize that the present invention could also be used to determine the least similarity between a profile and a target.

**[00142]** Temporal aspects of the present invention have not been largely discussed, however one of skill in the art will appreciate that computationally profiles are more time consuming than generation of transition probabilities. Thus, off-line or pre-generated profiles which are then compared against substantially real-time generated transition probabilities may be efficient use of compute power providing the user with a more rapid response than real-time generation of both profiles and transition probabilities.

**[00143]** As used in this description, "one embodiment" or "an embodiment" or similar phrases means that the feature(s) being described are included in at least one embodiment of the invention. References to "one embodiment" in this description do not necessarily refer to the same embodiment; however, neither are such embodiments mutually exclusive.

Nor does "one embodiment" imply that there is but a single embodiment of the invention. For example, a feature, structure, act, etc. described in "one embodiment" may also be included in other embodiments. Thus, the invention may include a variety of combinations and/or integrations of the embodiments described herein.

**[00144]** Thus a method and apparatus for fundamental operations on token sequences: computing similarity, extracting term values, and searching efficiently have been described.